

Testing the reconstruction of modelled particulate organic carbon from surface ecosystem components using PlankTOM12 and Machine Learning

Anna Denvil-Sommer¹, Erik T. Buitenhuis¹, Rainer Kiko^{2,3}, Fabien Lombard^{2,4}, Lionel Guidi², Corinne Le Quéré¹

¹School of Environmental Science, University of East Anglia, Norwich, UK

²Sorbonne Université, Centre National de la Recherche Scientifique (CNRS), Laboratoire d'Océanographie de Villefranche (LOV), Villefranche-sur-Mer, France

³GEOMAR Helmholtz Center for Ocean Research, Kiel, Germany

⁴Institut Universitaire de France (IUF), Paris, France

Correspondence to: Anna Denvil-Sommer (anna.sommer.lab@gmail.com)

Abstract. Understanding the relationship between surface marine ecosystems and the export of carbon to depth by sinking organic particles is key to represent the effect of ecosystem dynamics and diversity, and their evolution under multiple stressors, on the carbon cycle and climate in models. Recent observational technologies have greatly increased the amount of data available, both for the abundance of diverse plankton groups and for the concentration and properties of particulate organic carbon in the ocean interior. Here we use synthetic model data to test the potential of using Machine Learning (ML) to reproduce concentrations of particulate organic carbon within the ocean interior based on surface ecosystem and environmental data. We test two machine learning methods that differ in their approaches to data-fitting, the Random Forest and XGBoost methods. The synthetic data is sampled from the PlankTOM12 global biogeochemical model using the time and coordinates of existing observations. We test 27 different combinations of possible drivers to reconstruct small (POC_S) and large (POC_L) particulate organic carbon concentrations. We show that ML can successfully be used to reproduce modelled particulate organic carbon over most of the ocean based on ecosystem and modelled environmental drivers. XGBoost showed better results compared to Random Forest thanks to its gradient boosting trees architecture. The inclusion of Plankton Functional Types (PFTs) in driver sets improved the accuracy of the model reconstruction by 58% on average for POC_S, and by 22% for POC_L. Results were less robust over the Equatorial Pacific and some parts of the high latitudes. For POC_S reconstruction, the most important drivers were the depth level, temperature, microzooplankton and PO₄, while for POC_L it was the depth level, temperature, mixed-layer depth, microzooplankton, phaeocystis, PO₄ and chlorophyll *a* averaged over the mixed-layer depth. These results suggest that it will be possible to identify linkages between surface environmental and ecosystem structure and particulate organic carbon distribution within the ocean interior using real observations, and to use this knowledge to improve both our understanding of ecosystem dynamics and of their functional representation within models.

1. Introduction.

Progress in numerical ocean modelling over multiple decades coupled with fundamental knowledge of fluid dynamics have led to an explicit representation of ocean dynamics in Earth System Models and of most of its key features, apart from small-scale features which are parametrized. In contrast, ecosystem dynamics in ocean biogeochemical models are much more reliant on empirical data for growth and loss processes, with the theoretical basis limited to the dynamic representation of interactions among lower trophic levels (zooplankton and smaller organisms) and their influence on carbon pools and fluxes (Le Quéré et al., 2005; Hood et al., 2006). The recent advances in observational technologies including imaging data (Guidi et al., 2016), genomics (Kirchman et al., 2016), and field study (Mutshinda et al., 2017; Batten et al., 2019, Lombard et al., 2019), offer new opportunities to improve our understanding of marine ecosystem dynamics, and to better represent its influence on carbon pools and fluxes in models that are used to project future climate change and associated impacts on ecosystems.

One strategy to represent lower trophic interactions in global biogeochemical models is to combine different species into Plankton Functional Types (PFTs) based on their unique influence on global biogeochemical cycles (Le Quéré et al., 2005; Hood et al., 2006). This approach enables the representation of plankton types that are unique, have an

51 influence on other PFTs within the ecosystem and are of quantitative importance for carbon flux and other
52 biogeochemical fluxes. The PlankTOM12 model is among the most detailed in this category of models with its
53 inclusion of an explicit representation of twelve PFTs: six phytoplankton, five zooplankton, and bacteria.
54 PlankTOM12 builds on the published version PlankTOM10 (Le Quéré et al., 2016) that has been extended to include
55 gelatinous zooplankton (Wright et al., 2021) and pteropods (Buitenhuis et al., 2019). Much effort has been put into
56 the development of PFTs and associated representation of surface ecosystem dynamics, which has led to the
57 demonstration that: (1) the representation of trophic levels was a key determinant of the low chlorophyll *a*
58 concentration observed in the Southern Ocean summer (Le Quéré et al., 2016); (2) CaCO₃ dissolution above the
59 lysocline is needed to reproduce observations of both biomass and export of PFT calcifiers, and (3) gelatinous
60 zooplankton plays an important role in determining surface biomass of other PFTs (Wright et al., 2021).

61 In contrast, the transfer of organic matter resulting from surface ecosystem dynamics into carbon exported to the deep
62 ocean via the sinking of particulate organic matter has received much less attention, so that improvements in the
63 representation of the PFTs do not necessarily translate into improvements in sinking of particulate matter (Wright et
64 al., 2021). The export flux of particulate organic carbon from the surface ocean to depth is around 10 PgC yr⁻¹
65 (Schlitzer, 2002), which is as large as the CO₂ emitted to the atmosphere by human activities and nearly four times
66 larger than the mean oceanic CO₂ sink in recent decades (Friedlingstein et al., 2022). Changes in carbon exported to
67 depth can have a large impact on air-sea CO₂ fluxes and on the amount of CO₂ emissions that remain in the atmosphere
68 where they cause climate change.

69 The growing amount of observations provides the opportunity to develop a new approach to explore the linkages
70 between surface ecosystem dynamics and the distribution of particulate organic carbon in the ocean, and to improve
71 the representation of particle sinking fluxes in models. However, there is a risk of over-interpreting the data by
72 applying Machine Learning (ML) methods directly to link the observed surface environment and ecosystem structure
73 with the observed particulate organic carbon distribution. The use of synthetic observations based on model data
74 therefore provides a minimum test to assess the likely success and usefulness of such an approach.

75 ML has been widely used in biogeochemical and geophysical applications and provided efficient results in
76 reconstructions of ocean surface pCO₂ (Friedrich and Oschlies, 2009; Telszewski et al., 2009; Landschützer et al.,
77 2013; Denvil-Sommer et al., 2019) and of particulate organic carbon (Sauzède et al., 2016, 2017) as well as in the
78 analysis of driver importance (Sauzède et al., 2020).

79 Here we use model data to verify the hypothesis that the composition of surface ecosystems and environmental
80 conditions are indeed reflected in the abundance and size of the organic particles in the ocean interior. We reconstruct
81 the concentration of organic particles as represented by small (POC_S, particles < 256µm) and large (POC_L, particles >
82 256µm) particulate carbon in the PlankTOM12 model. Using this information alongside with modelled environmental
83 and ecosystem conditions we develop a ML method to reproduce POC_S and POC_L over the global ocean and verify
84 the hypothesis. This constitutes a necessary although not sufficient test that the approach can subsequently be used to
85 reveal linkages using real observations and to inform model developments.

86 2. Data and Methods.

87 In this section we describe a set of variables that will be used to test the ML method's ability to reconstruct particulate
88 organic carbon concentrations based on ocean model data. We create a set of synthetic data by sampling a model at
89 the time and location of real-world observations. We discuss the availability and distribution of real-world
90 observations and their limitations. In this section we also describe the PlankTOM12 global ocean biogeochemical
91 model and how we use it to develop a ML method and test its ability to reconstruct small and large particulate organic
92 carbon with a limited number of observations. To provide resemblance to the real data availability we focus on the
93 period 2009-2013 which guarantees additional sampling of co-located biological, chemical, and environmental
94 variables from the Tara expeditions (Sunagawa et al., 2020).

95 Two sets of data are needed to test the Machine Learning method: a set of targets and a set of drivers. The drivers
96 represent the input variables to the ML method (here the biological, chemical, and environmental variables). The
97 targets represent the variables we are trying to reconstruct (here the particulate organic matter POC_S and POC_L). The
98 ML will then determine the relationship between the drivers and targets, which can then be applied in regions where
99 drivers are available to infer targets where the later data do not exist.

100

101

2.1.1. Measurements of particle size distributions and concentrations (the targets).

102

We use observations of particle distribution in two ways. First to determine the time and location of the observations, and second to verify that the ocean model is of sufficient quality to be used in this analysis. The sampling of the particulate organic carbon concentration is based on the data from an Underwater Vision Profiler 5 (UVP5) (Gorsky et al., 2000, 1992; Picheral et al., 2010; Kiko et al., 2022). UVP5 measures particles of size from 50 μm to a few mm. For the purpose of comparing the UVP5 data with the PlankTOM12 model data, we converted measured biovolume concentration (mm^3/L) of particles to carbon biomass concentrations ($\mu\text{mol}/\text{L}$) using the empirical equation from Alldredge (1998) for particulate organic carbon:

103

104

105

106

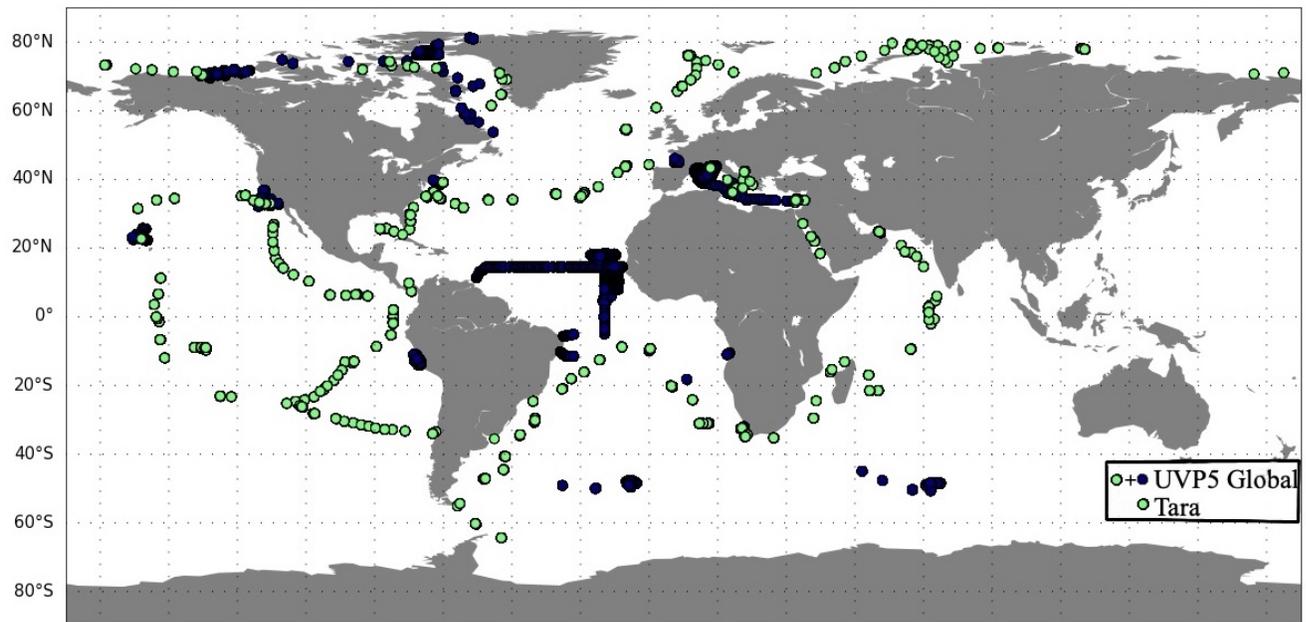
107

108

$$\text{BM} (\mu\text{g}) = 0.99 * \text{BV} (\text{mm}^3)^{0.52}$$

109

110



111

112

113

Figure 1. Location of the observations from the UVP5 database over the period 2009-2013. Green dots correspond to Tara expeditions, and were included in the global UVP5 database.

114

We summed size classes from 50.8 μm to 256 μm for the small particulate organic carbon (POC_S) and from 256 μm to 5.16 mm for large particulate organic carbon (POC_L). POCs below 100 μm are not well captured by the UVP sensor, which therefore underestimates this size-class of aggregated particles. We extrapolated the total size of particles up to 0.001mm by using the size spectra theory to provide a better estimate of POC biomass concentration in line with the model. Following Guidi et al. (2008) we used the abundance of particles sized from 0.250 to 1.5 mm excluding rare particles to estimate the coefficients of logarithmic relationship between the size of particles and its abundance:

115

116

117

118

119

$$\log(\text{abundance}) = a * \log(\text{size}) + b$$

120

121

Using this equation we estimated the abundance of particles of size less than 100 μm .

122

There are 2603 vertical profiles of UVP5 measurements during 2009-2013, including 752 profiles which are collocated with the stations from the Tara expeditions that provide the environmental and ecosystem variables (Figure 1; Section 2.1.2). The measurements are sparse in time and space. There are no measurements in the Southern Ocean, Western Pacific Ocean and Eastern Indian Oceans.

123

124

125

126

2.1.2. Measurements of environmental and ecosystem variables (the drivers).

127

We use observations of environmental and ecosystem variables to determine the time and location of the observations that are collocated with the target variables. To represent the main physical and chemical drivers responsible for the concentration and variability of POC_S and POC_L we use measurements of ocean temperature, chlorophyll *a*, phosphate PO₄, nitrates NO₃, mixed-layer depth (MLD). These variables were measured during Tara expeditions along with the

128

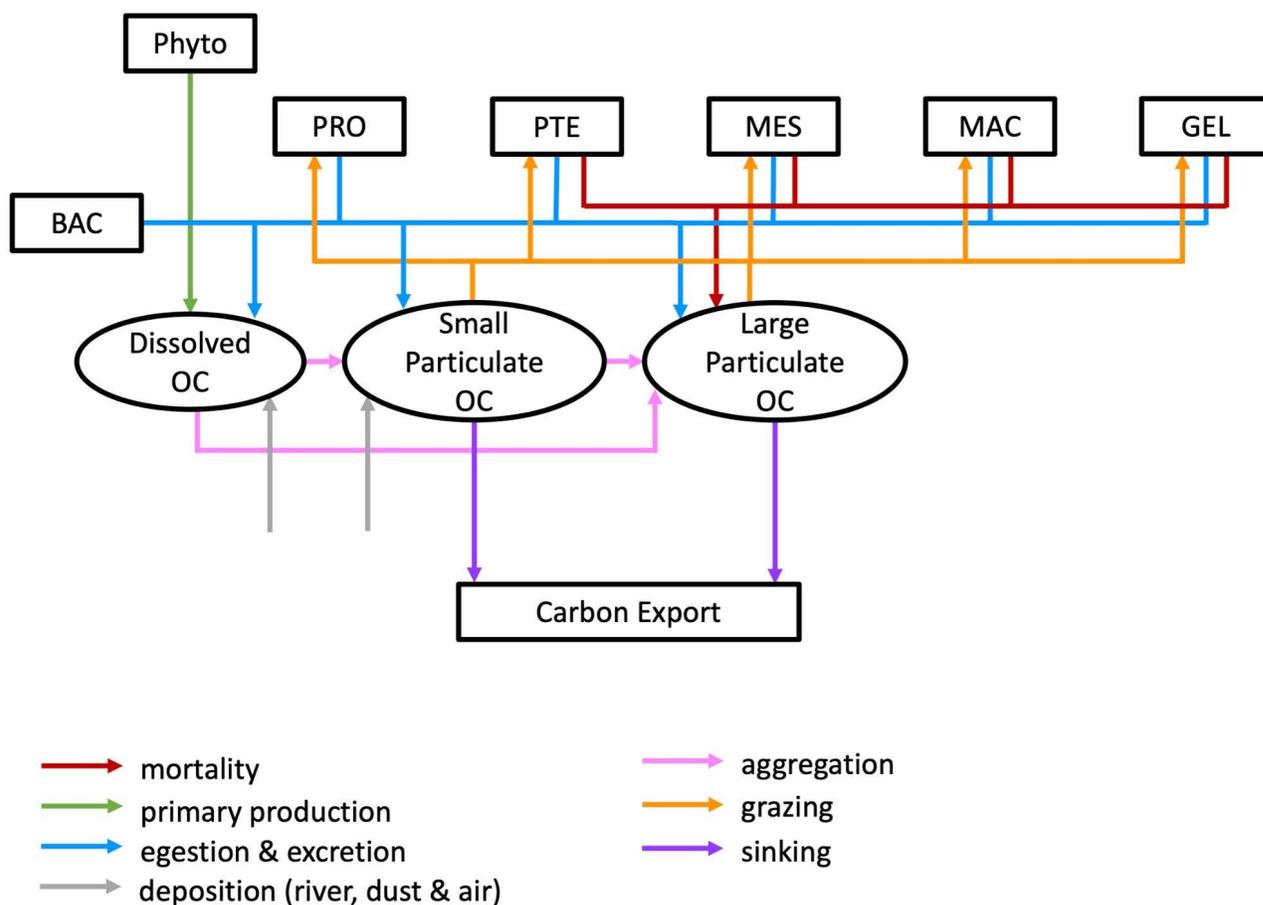
129

130

131 particle size distributions and concentrations using UVP instruments onboard these cruises. However, chlorophyll *a*,
 132 PO₄, and nitrates were not measured systematically at each depth level. Thus, their averages over MLD are tested as
 133 possible drivers as well. To represent the biological drivers, we use information on PFTs.

134 **2.1.3. The NEMO-PlankTOM12 Global biogeochemical model.**

135 We used the output from the NEMO-PlankTOM12 coupled physical–biogeochemical model of the global ocean at
 136 daily and monthly time resolution. NEMO represents physical transport processes and is used in its v3.6-ORCA2
 137 version, with a horizontal resolution 2° longitude and 0.3° to 1.5° latitude, and 31 vertical levels. It is forced by daily
 138 meteorological data from NCEP reanalysis (Kalnay et al., 1996) over the period 1948-2020, with output for 2009-
 139 2013 used here. This model version is identical to that used to estimate the ocean CO₂ sink in the Global Carbon
 140 Budget 2021 annual update (Friedlingstein et al. 2021).



141
 142 **Figure 2. Schematic representation of the flow of matter in and out of the two particulate organic carbon (OC) components**
 143 **of the PlankTOM12 marine ecosystem model. The various boxes represent: Phyto - phytoplankton that includes diatoms**
 144 **(DIA), mixed phytoplankton (MIX), coccolithophore (COC), picophytoplankton (PIC), phaeocystis (PHA) and N₂-fixers**
 145 **(FIX); PRO - protozooplankton, PTE - pteropod, MES - mesozooplankton, MAC - macrozooplankton, GEL - gelatinous**
 146 **zooplankton, BAC - bacteria.**

147 PlankTOM12 represents ecosystem dynamics based on the representation of 12 PFTs: diatoms (DIA), mixed
 148 phytoplankton (MIX), coccolithophore (COC), picophytoplankton (PIC), phaeocystis (PHA), N₂-fixers (FIX), micro-
 149 or protozooplankton (PRO), pteropod (PTE), mesozooplankton (MES), gelatinous zooplankton (GEL), and bacteria
 150 (BAC). PlankTOM12 keeps track of the carbon biomass (μmol/L) of these PFTs over model depth levels resulting
 151 from environmental and ecosystem processes and their interactions (Le Quéré et al. 2016).

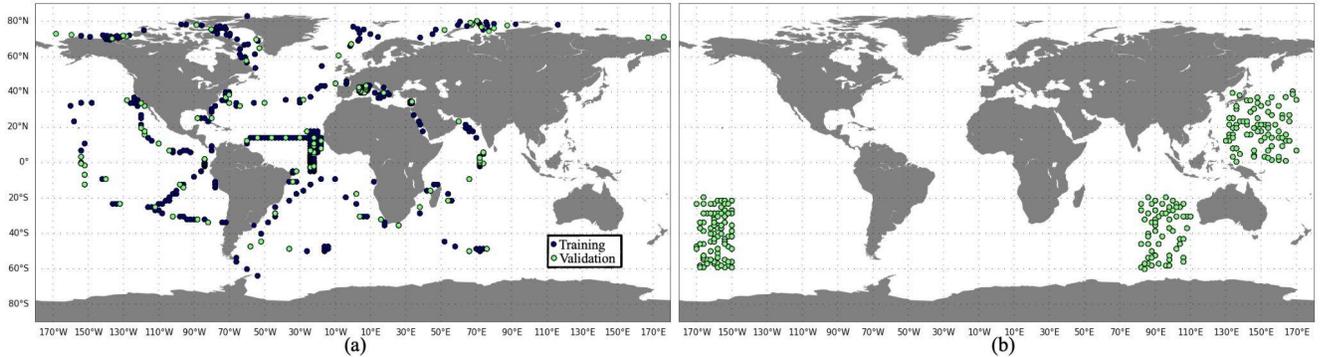
152 PlankTOM12 represents sinking processes through the explicit representation of two organic particle of different size,
153 with small particles sinking at a constant speed of 3 m/d, and larger particles sinking at a variable speed between 3
154 and 150 m/d depending on the ballast effect of their mineral content (Buitenhuis et al., 2013). In addition, a dissolved
155 organic carbon component is transported via ocean currents. Particles are generated through mass flux from the PFTs
156 resulting from mortality and egestion and from aggregation through differential sinking or turbulent coagulation, and
157 destroyed through grazing by zooplankton and remineralisation by bacteria and through disaggregation from shear
158 currents. Large PFTs contribute mostly to POC_L, while small PFTs contribute mostly to POCs. (Le Quéré et al. 2016;
159 Fig. 2).

160 The NEMO-PlankTOM12 model output was sampled at the time and location identified from the observations
161 mentioned above to create a synthetic dataset. The model grid-coordinate closest to the real geographical position was
162 chosen. If several measurements were co-localised at the same grid coordinate and same time step (day for daily
163 PlankTOM12 and month for monthly PlankTOM12 outputs), it is counted as one measurement. This model sampling
164 produced 400 positions when using the daily or monthly PlankTOM12 outputs. All drivers and targets were taken
165 from the model output at the corresponding coordinates up to 1400 m depth. These outputs served as the reference for
166 validation and evaluation of the ML methods and for establishing the sets of the most important drivers.

167 **2.2. Method.**

168 We tested 2 ML methods that are widely used in target's reconstruction based on tabular data sets: the Random Forest
169 regressor and the XGBoost (Extreme Gradient Boosting) regressor. The Random Forest (RF) regressor is an ensemble
170 algorithm that contains a number of decision trees on various subsets of the given dataset and takes as output the
171 average of prediction from each tree estimator. RF can run several trees at the same time allowing a use of a large
172 number of input variables, and it is robust to overfitting (Biau, 2012). XGBoost (XGB) regressor is an effective tree-
173 based ensemble learning algorithm (Chen and Guestrin, 2016). It builds several models sequentially where each new
174 model attempts to correct errors from the previous one. XGBoost uses the gradient descent algorithm to minimise the
175 loss function of the model. Using RF and XGBoost we can estimate the driver importance to identify which driver has
176 the greatest impact on the predictions. To check the driver importance, we use `drop_col_feat_imp` python function
177 (<https://gist.github.com/erykml/6854134220276b1a50862aa486a44192>). This method estimates how the accuracy of
178 the ML output changes if one of the drivers is dropped off from a driver set (DS) based on the training dataset.

179 Effective ML algorithm requires sets of training, validation and test data. The training data builds up the ML model.
180 Model evaluates training data repeatedly to learn about the relationship between inputs (driver set) and known outputs
181 (target set) and adjusts itself to better represent the target. The purpose of validation data is to evaluate the model
182 during its training by introducing new unseen data. It allows us to evaluate how a developed model works on a new
183 dataset and to optimise hyperparameters. The test data evaluate the final accuracy of the ML model and confirm that
184 the model works correctly on any unseen data. It is new data that did not participate in the training algorithm. The
185 accuracy is worse on validation and test data compared to training data set. The difference in model performance on
186 training and validation data can signal an overfitting, while this difference between validation and test data can
187 demonstrate an effect of data mismatch. It is worth noting that RF does not necessarily need validation data set as they
188 perform internal validation. During the training algorithm each tree is constructed from a random subset of original
189 data, usually it represents two thirds of data and one third of data is used to estimate out-of bag error to assess model
190 performance. XGB uses a validation data set to evaluate the model during training and to prevent overfitting by
191 applying an early stopping. In the present study the available data were split into training and validation data sets (Fig.
192 3a). Validation data is not included in RF training, however we use it to test the performance of trained RF and tune
193 hyperparameters afterwards. The test data are taken from the regions where there are no observations (Fig. 3b): 3
194 months for each year from the period 2009-2013 and 6 positions for each month were chosen randomly. This will
195 allow us to identify the possible accuracy of reconstruction that can be reached in these regions when we will apply a
196 developed method to real observations. However, when POCs and POC_L will be reconstructed using only real-world
197 observations, we will need to split all available data into training, validation and test data sets.



198
199 **Figure 3. The spatial distribution of: (a) - training (blue) and validation (green) data sets; (b) - test data set; based on**
200 **PlankTOM12 monthly outputs.**

201 We use RandomForestRegressor function from scikit-learn (<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html>) with its default parameters and
202 min_sample_leaf equals 20. To apply XGBoost regressor we use XGBRegressor from xgboost
203 (https://xgboost.readthedocs.io/en/stable/python/python_intro.html). Parameters were set as follows
204 n_estimators=2000, max_depth=7, eta=0.01, subsample=0.7, colsample_bytree=0.8, gamma=0.01 for POC_L and
205 gamma= 0.3 for POC_s, early_stopping_rounds = 10.
206

207 We tested 27 driver sets (DSs) that are summarised in Table 1. For each DS we identify the most important drivers
208 that influenced the reconstruction of small (POC_s) and large (POC_L) particulate organic carbon concentration. The
209 drivers include geographic variables (depth, sin(latitude), cos(longitude)), physical variables (incident light, MLD,
210 co-located temperature), chemical variables (PO₄, NO₃, including co-located values and averages over the MLD), and
211 biological variables (chlorophyll *a*, 12 PFTs listed above: DIA, MIX, COC, PIC, PHA, FIX, PRO, PTE, MES, GEL,
212 BAC, including co-located values and averages over the MLD).

213 **Table 1. Compounds of driver's sets: dark grey cells correspond to the drivers present in the driver set. 'vp' – vertical**
214 **profile, 'mean' – average over MLD, 'back' – values from previous month.**

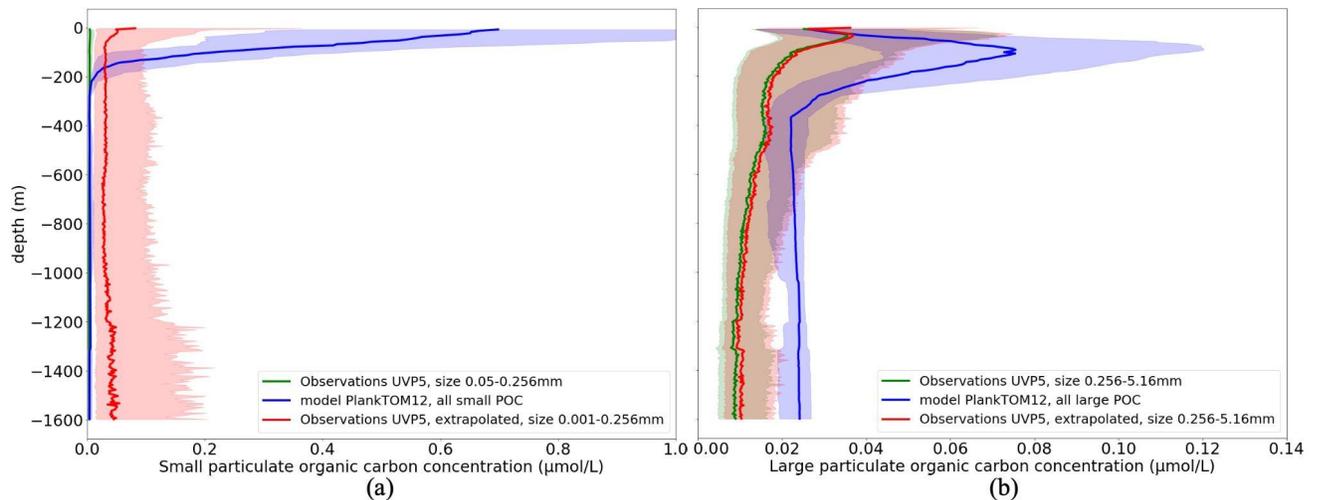
Driver set	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27
depth																											
Sin(lat)																											
Sin(long)																											
Cos(long)																											
Incident light																											
MLD																											
Temperature vp																											
CHL vp																											
CHL mean																											
NO ₃ vp																											
PO ₄ vp																											
NO ₃ mean																											
PO ₄ mean																											
BAC vp																											
MES vp																											
PTE vp																											
DIA vp																											
COC vp																											
PIC vp																											
PHA vp																											
GEL vp																											
PRO vp																											
MAC vp																											
MIX vp																											
FIX vp																											
BAC mean																											
MES mean																											
PTE mean																											
DIA mean																											

257

3.1. Data analysis.

258 In this study we test the capacity to reconstruct particulate organic carbon from sparse observations by using ML and
259 a synthetic data set based on the PlankTOM12 model output. We compare observations and the output of the ocean
260 model to provide a minimum of validation for the model data and to help explain differences in ML results when
261 applied to real observations in the future.

262 Figure 4 shows the vertical profile of small (POC_S) (Fig.4a) and large (POC_L) particulate organic carbon (Fig.4b)
263 based on the median from observations (green) and from daily PlankTOM12 model output (blue). Shading
264 corresponds to values between the 0.25 and 0.75 percentiles.



265 **Figure 4. Comparison of the vertical distribution of particulate organic carbon concentrations (µmol/l) from UVP5**
266 **measurements (green), PlankTOM12 daily model (blue) and extrapolated UVP5 measurements (red): (a) - small particulate**
267 **organic carbon concentrations; (b) - large particulate organic carbon concentrations. The median is shown in dark and the**
268 **shading corresponds to values between the 0.25 and 0.75 percentiles. The size of the particles does not correspond**
269 **completely between the observations and the model, for POC_L the UVP particle range is chosen as 0.256-5.16 mm that**
270 **corresponds approximately to the POC_L in the model.**
271

272 PlankTOM12 overestimates POC_S up to 3 µmol/L in the first 200m (Fig.4a, green and blue curves). UVP5 does not
273 capture all small particles that is why we extrapolated the size range of UVP measurements (red curve, see details in
274 2.1.1). The extrapolated measurements show an increase in POC_S in the first 100m, however this increase still results
275 in the lower concentration compared with PlankTOM12. These results indicate that PlankTOM12 overestimates the
276 concentration of small particulate organic carbon. PlankTOM12 also overestimates POC_L by up to 0.08 µmol/L
277 in the first 200m and does not catch the increase in POC_L between 300 and 500m. Observations show an increase in
278 POC_L concentration in the first 50m while PlankTOM12 reproduces it lower, at 100m. The RMSE between modelled
279 and observed POC_S is 0.33 µmol/L, with correlation coefficient equals 0.083. RMSE equals 0.23 µmol/L with
280 correlation coefficient 0.061 for POC_L. The exclusion of isolated large values of POC_L (>2 µmol/L) from the
281 observation data set reduces the RMSE of POC_L to 0.062 µmol/L with correlation 0.18. We believe that these
282 differences result from differences in space and time resolution of observations and ocean model outputs. *In-situ*
283 measurements are obtained at a particular time of the day and a particular latitude-longitude position while the model
284 provides estimations over the day (or month) and on the model grid (2° longitude and mean 1.1° latitude resolution).

285 We concluded that observed and modelled POC_S and POC_L have a common tendency in their vertical distributions.
286 However, among other things, differences in amplitudes may affect our findings in this work when we develop a ML
287 method based on observations only.

288 Due to the constraint in data availability further we use monthly PlankTOM12.

289 Before developing a ML method, we investigate the interactions between targets and drivers in the model. Table 2
290 shows the correlation coefficients between the POC_S and POC_L and corresponding drivers that can influence POC_S
291 and POC_L variability. **Correlation between drivers could also provide valuable information to minimise the number of**

292 driver but they are not shown here where the focus is on discovering the effect of a large set of drivers on POC
 293 distribution, and because driver correlations could also result from the physics as well as from the model construction.
 294 POC_s correlates with gelatinous zooplankton (GEL, $r=0.66$), microzooplankton (PRO, $r=0.63$), coccolithophore
 295 (COC, $r=0.56$), as well as with their values from previous time step (GEL, $r=0.67$; PRO, $r=0.51$; COC, $r=0.59$).
 296 Coccolithophore is one of the most abundant phytoplankton types in this version of the PlankTOM model (similar to
 297 Wright et al., 2021). The growth of phytoplankton transfers dissolved inorganic carbon into dissolved organic carbon
 298 which further aggregates into POC_s and POC_L. Also, POC_s is generated from microzooplankton egestion and
 299 excretion (Fig. 2). In addition to the mentioned above PFTs, POC_s shows a correlation 0.44 with temperature vertical
 300 profile at both the considered time step and at the previous time step. POC_s has a negative correlation with NO₃ ($r=-$
 301 0.46) and PO₄ ($r=-0.41$).

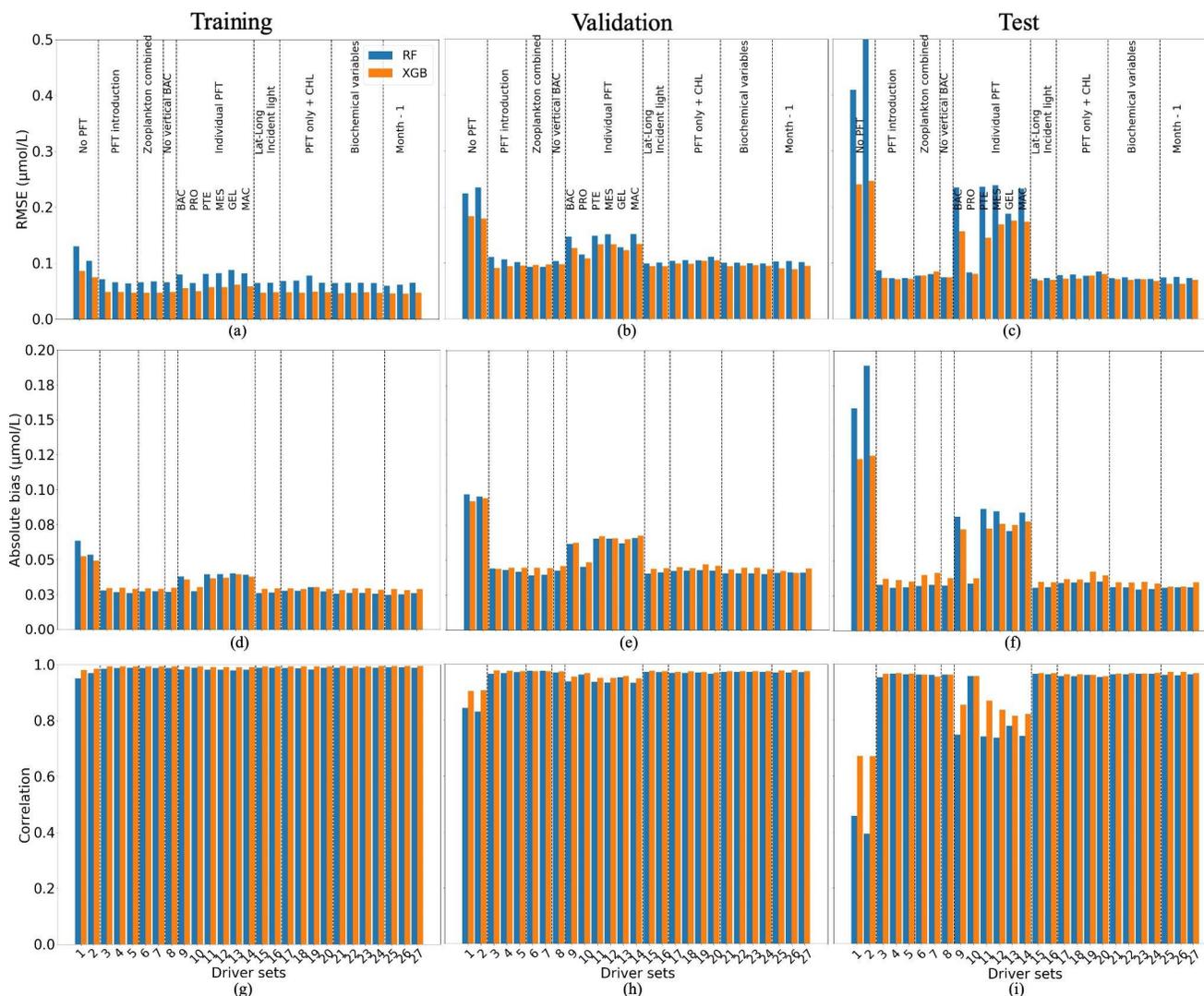
302 POC_L does not show a high correlation with any of the proposed drivers individually and is therefore most likely the
 303 result of multiple processes and/or multiple drivers, including for its production and destruction. The ML approach
 304 should be able to identify combinations of drivers beyond straight correlations that are investigated directly here.
 305 POC_L has the highest correlation with chlorophyll *a* ($r=0.42$), gelatinous zooplankton at the considered time step
 306 ($r=0.37$), and at previous time step ($r=0.36$). Gelatinous zooplankton contribute to POC_L formation through egestion
 307 and excretion mainly from mucus (Fig. 2). As explained in Wright et al. (2021), mucus forms a large low-density mass
 308 through aggregation with other particles. It can explain a correlation of gelatinous zooplankton with POC_L in
 309 PlankTOM12.

310 **Table 2. Correlation coefficient between small (POC_s) and large (POC_L) particulate organic carbon concentration and**
 311 **possible drivers. Estimation is based on monthly PlankTOM12 output at the position of real-world observations from Fig.**
 312 **1. ‘vp’ – vertical profile, ‘mean’ – average over MLD, ‘back’ – values from previous month.**
 313

Driver	POC _s	POC _L	Driver	POC _s	POC _L	Driver	POC _s	POC _L
POC	1.00	0.33	BAC vp	-0.14	0.15	BAC back vp	-0.10	0.09
GOC	0.33	1.00	MES vp	-0.09	0.07	MES back vp	-0.09	-0.07
Depth	-0.32	-0.24	PTE vp	-0.07	0.17	PTE back vp	-0.08	0.08
Temperature vp	0.44	0.17	DIA vp	-0.04	0.15	DIA back vp	-0.03	0.09
Temp back vp	0.44	0.17	COC vp	0.56	0.31	COC back vp	0.60	0.31
MLD	-0.01	-0.07	PIC vp	0.00	0.07	PIC back vp	0.06	0.06
NO ₃ vp	-0.46	0.01	PHA vp	0.27	0.15	PHA back vp	0.30	0.17
PO ₄ vp	-0.41	0.04	GEL vp	0.66	0.37	GEL back vp	0.68	0.36
NO ₃ back vp	-0.46	0.03	PRO vp	0.63	0.16	PRO back vp	0.51	0.14
PO ₄ back vp	-0.41	0.05	MAC vp	0.07	0.14	MAC back vp	0.08	0.13
CHL vp	0.18	0.42	MIX vp	0.07	0.17	MIX back vp	0.03	0.05
CHL back vp	0.11	0.22	FIX vp	-0.00	0.23	FIX back vp	-0.00	0.23

314

315 **3.2. Development of the Machine Learning method.**



316

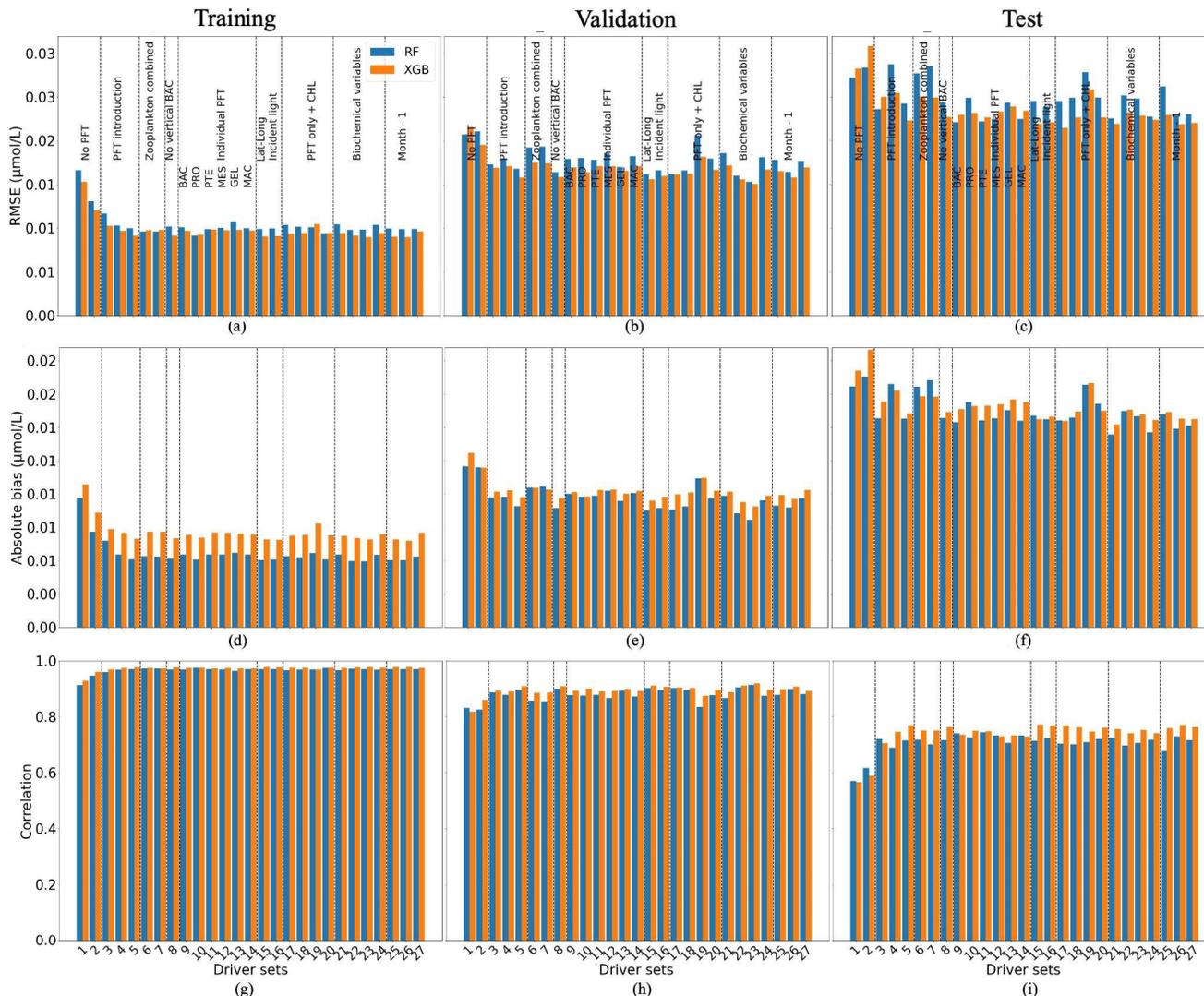
317 **Figure 5. Comparison of the performance of the Random Forest (RF) and XGBoost methods and their fit to data for small**
 318 **(POC_s) particulate organic carbon concentration; (a, b, c) - RMSE in µmol/l, (d, e, f) - absolute bias in µmol/, (g, h, i) -**
 319 **correlation coefficient; (a, d, g) - training data set, (b, e, h) - the validation data set, (e, f, i) - the test data set. Results**
 320 **compare data from the original (sampled) PlankTOM12 model output and POC_s reconstructed using RF (blue) and XGB**
 321 **(orange). The low RMSE and absolute biases indicate better performance of the ML method.**

322 We tested 27 sets of drivers (Table 1) and two ML methods, Random Forest (RF) and XGBoost regression (XGB).

323 Figure 5 shows the statistics of POC_s reconstruction using RF and XGB. XGB (orange) generally overperforms RF
 324 (blue). The statistics are slightly worse for the validation and test data sets, as expected. For reconstructions using
 325 XGB, the RMSE and absolute bias are about 0.05 µmol/L and 0.03 µmol/L on the training data set and vary around
 326 0.1 µmol/L and 0.05 µmol/L, on the validation and test data, respectively. Correlation coefficients (Fig. 5g, h, i) have
 327 high values on all datasets showing that the vertical profiles of POC_s have a correct shape. These results show that the
 328 available spatial and temporal coverage of *in situ* observations can be sufficient to reconstruct POC_s with an
 329 appropriate accuracy over the global ocean. The analysis of global maps (shown below) will help to identify areas
 330 with low accuracy and their differences with training regions.

331 The worse results (highest RMSE, highest absolute bias, lowest correlation) are produced when there are no PFTs in
 332 the driver set (DS1 and DS2; Figure 5): for XBoost, RMSEs are 0.24 µmol/L, absolute biases equal to 0.12 µmol/L
 333 with correlation coefficient 0.67 on the test data sets. Poor results are also obtained for DS9, 11, 12, 13 and 14: these
 334 5 driver sets do not have any information on microzooplankton (PRO) and show high RMSEs and absolute biases,

335 around 0.16 $\mu\text{mol/L}$ and 0.074 $\mu\text{mol/L}$, with low correlation, 0.83, compared with other driver sets which include
 336 PRO. These results indicate that microzooplankton plays an important role in POC_s variability in the PlankTOM12
 337 model.



338
 339 **Figure 6. Comparison of the performance of the Random Forest (RF) and XGBoost methods and their fit to data for large**
 340 **(POC_L) particulate organic carbon concentration; (a, b, c) - RMSE in $\mu\text{mol/L}$, (d, e, f) - absolute bias in $\mu\text{mol/L}$, (g, h, i) -**
 341 **correlation coefficient; (a, d, g) - training data set, (b, e, h) - the validation data set, (e, f, i) - the test data set. Results**
 342 **compare data from the original (sampled) PlankTOM12 model output and POC_L reconstructed using RF (blue) and XGB**
 343 **(orange). The low RMSE and absolute biases indicate better performance of the ML method.**

344 Figure 6 shows the statistics of POC_L reconstruction using RF and XGB. XGBoost again slightly overperforms RF on
 345 most driver sets. Results for driver sets with PFTs show lower RMSEs and absolute biases, and higher correlation
 346 coefficients. Except for the effect of PFTs on the POC_L reconstruction, we did not observe a clear influence of one
 347 driver or group of drivers. Using XGBoost the reconstruction of POC_L shows the RMSE in DS1 is high at 0.03 $\mu\text{mol/L}$,
 348 while it is in the range of 0.021-0.026 $\mu\text{mol/L}$ in DS3-DS27, with absolute bias in DS1 of 0.02 $\mu\text{mol/L}$ and 0.015-
 349 0.018 $\mu\text{mol/L}$ for DS3-DS27 based on test data (Fig. 6c, f). Likewise, a correlation coefficient of 0.56 for DS1, and
 350 between 0.7 and 0.77 for DS3-DS27 based on the training data set (Fig. 6g).

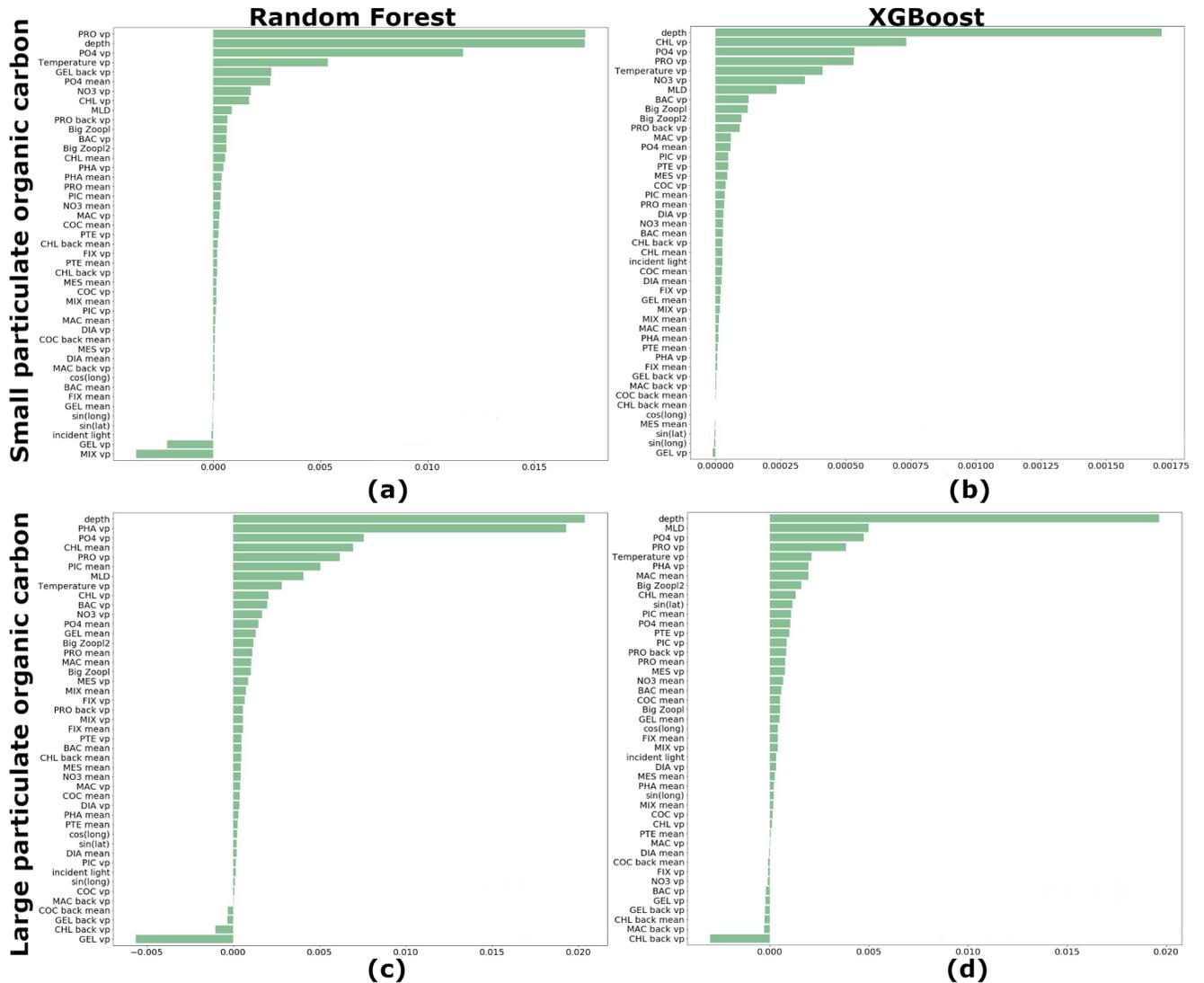
351 We estimated the ranking of importance for each driver averaged over 27 driver sets (Table 1) for RF and XGB (Fig.
 352 7). Both, RF and XGB, show that microzooplankton (PRO), depth level, temperature, NO₃ and PO₄ play a dominant
 353 role in reconstruction of POC_s. The absence of gelatinous zooplankton (GEL) can slightly improve the reconstruction.

354 Also, latitude and longitude do not affect POC_S reconstruction. The depth level, temperature, MLD, microzooplankton
355 (PRO) and phaeocystis (PHA), PO₄, and chlorophyll *a* averaged over MLD play a dominant role in POC_L
356 reconstruction.

357 The sinus of latitude is in the top ten drivers that most affect POC_L using XGBoost method: POC_L distribution has a
358 lot of meridional variability that results in the sinus of latitude being in the top 10 drivers. As for POC_S, gelatinous
359 zooplankton (GEL) shows a negative rank of driver importance and its removal from the list of drivers can improve
360 the statistics of reconstruction. Also, chlorophyll *a* concentration from the previous month shows a similar effect on
361 POC_L (Fig. 7c, d).

362 It is worth noting that any driver that shows negative importance in the reconstruction has only a small influence on
363 the accuracy (Fig. 5 and 6). Thus, its removal does not improve the reconstruction significantly.

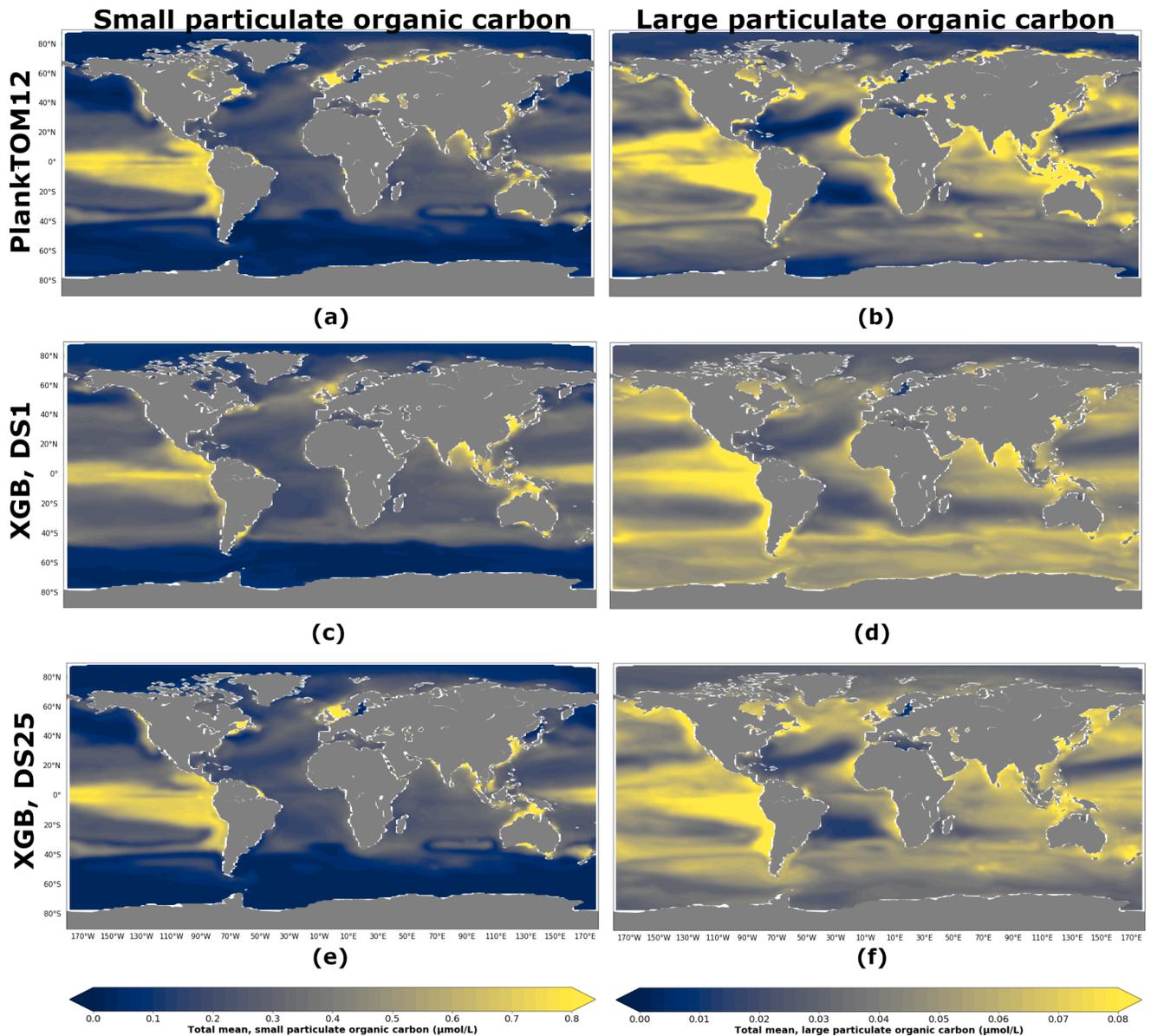
364 Based on Figures 5, 6 and 7 we have chosen 10 driver sets with low RMSEs and absolute biases, and high correlation
365 coefficients (based on test data set) for POC_S and POC_L to provide global maps of these statistics and to see their
366 regional distributions. DS 5, 15, 16, 21, 22, 23, 24, 25, 26, 27 were chosen for further investigation of POC_S
367 reconstruction; DS 5, 8, 15, 16, 17, 21, 23, 25, 26, 27 – for POC_L reconstruction. Common for POC_S and POC_L driver
368 sets 5, 15, 16, 21, 23, 25, 26, 27 include all PFTs and their average over MLD, geographical positions and incident
369 light as well as chlorophyll *a*, PO₄, and gelatinous zooplankton and microzooplankton from the previous time step
370 (Table 1). Also, we found that POC_S reconstructions rest on biochemical conditions (DSs 21 and 24), while POC_L
371 reconstruction mostly depends on the composition of the PFTs in the driver set (DSs 8 and 17). Additionally, we keep
372 DS1 to demonstrate a global effect of PFTs on reconstruction.



373
 374 **Figure 7. Ranking of importance for each driver averaged over 27 driver sets: (a) - Random Forest (RF) for reconstruction**
 375 **of small (POC_S) particulate organic carbon concentration; (b) - XGBoost (XGB) for small (POC_S) particulate organic**
 376 **carbon concentration; (c) - RF for POC_L concentration; (d) - XGB for POC_L concentration. ‘vp’ – vertical profile, ‘mean’**
 377 **– average over MLD, ‘back’ – values from previous month.**

378 **3.3. POC_S and POC_L vertical profile reconstruction over the global ocean**

379 In the previous section we showed that XGBoost provides the best results for the reconstructions of POC_S and POC_L.
 380 Further we use this ML method. Here we will discuss the regional results of DS1 without PFTs and 10 best driver sets
 381 chosen for each target separately.



382

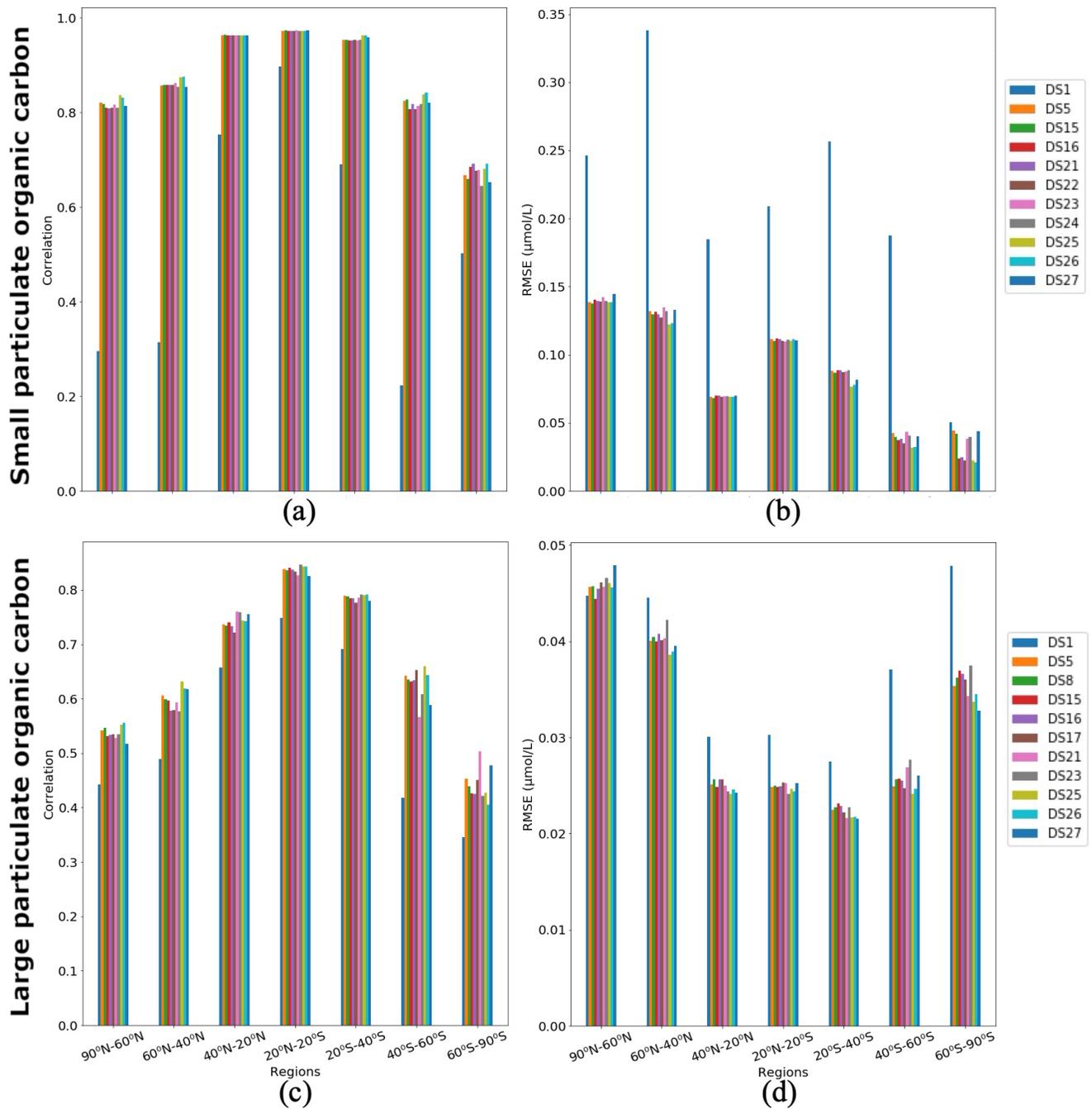
383 **Figure 8. Total averaged over the depth and period 2009-2013 small (POC_S) and large (POC_L) particulate organic carbon**
 384 **concentration: (a) – PlankTOM12 POC_S, (b) – PlankTOM12 POC_L, (c) – reconstruction of POC_S based on DS1 (NoPFT)**
 385 **using XGBoost, (d) – reconstruction of POC_L based on DS1 using XGBoos, (e) - reconstruction of POC_S based on DS25**
 386 **(vertical profiles of zooplanktons, and zooplankton and phytoplankton averaged over MLD) using XGBoost, (f) -**
 387 **reconstruction of POC_L based on DS25 using XGBoost.**

388 Figure 8 shows POC_S and POC_L concentration averaged over the depth and period 2009-2013 for PlankTOM12 (Fig.
 389 8a, b), XGboost reconstruction based on DS1 (Fig. 8c, d) and XGBoost reconstruction based on DS25 (Fig.8 e, f).
 390 XGBoost captures well the spatial patterns: the high concentration of POC_S in the Equatorial Eastern Pacific and its
 391 low concentration at high latitudes, as well as the high concentration of POC_L in the Equatorial Eastern Pacific and in
 392 the North of the Indian Ocean and its low concentration in the Subtropical North and South Atlantic and in the
 393 Subtropical North Pacific. The presence of PFTs in driver sets (Fig. 8e, f) improves the reconstruction: the spatial
 394 patterns and its amplitude are visually close to ones from PlankTOM12 (Fig. 8a, b). The high concentration of POC_S
 395 in the Equatorial Eastern Pacific is represented better using DS25 compared with DS1 where the concentration in the
 396 latitude band 0°S-20°S along the Peru is overestimated. Also, small decreases of POC_S in the Subtropical North and
 397 South Atlantic are captured better when we use DS25. Similar for POC_S, the high concentration in the Equatorial
 398 Eastern Pacific is represented better using DS25 compared with DS1 where the concentration misses the small

399 decrease between 20°N and 0°N. Also, small decreases of POC_L in the Subtropical North and South Atlantic as well
400 as in the Subtropical North Pacific are pronounced better with DS25.

401 Figure 9 shows regional correlation coefficients and RMSEs between PlankTOM12 and XGBoost reconstruction over
402 the global ocean for 2009-2013. We averaged correlation coefficient and RMSEs over 7 latitude zones: 90°N-60°N,
403 60°N-40°N, 40°N-20°N, 20°N-20°S, 20°S-40°S, 40°S-60°S, 60°S-90°S. In POC_s reconstruction, the DS1 shows the
404 lowest correlation across latitude bands (between 0.22 and 0.9), and highest RMSEs (0.05-0.34 μmol/L; Fig.9a, b).
405 DSs 25 and 26 show the highest correlations in the range of 0.68 (in region 60°S-90°S) and 0.97 (in region 20°N-20°S)
406 and the lowest RMSEs in the range of 0.021 (in region 60°S-90°S) and 0.14 μmol/L (in region 90°N-60°N). DS25
407 contains information on the previous-month distribution for micro-, macrozooplankton and gelatinous zooplankton
408 vertical profiles as well as coccolithophores and chlorophyll *a* averaged over the MLD. DS26 is like DS25 but the
409 drivers which bring information from the previous month are microzooplankton and gelatinous zooplankton vertical
410 profiles.

411 10 driver sets (excluding DS1) show their highest RMSEs in POC_s reconstruction in the region 90°N-60°N, with
412 values up to 0.14 μmol/L in DS27 (Fig. 9b). Figure 10 shows maps of RMSEs (a, b) and correlation coefficients (c,
413 d) between PlankTOM12 and reconstructed small particulate organic carbon (POC_s) by XGBoost using driver sets 1
414 (a, c) and 25 (b, d). The region 90°N-60°N shows improvement in RMSEs and absolute biases in DS25 compared with
415 DS1, with RMSEs decreasing from 0.2 μmol/L to 0.03 μmol/L in Norwegian Sea, Baffin Bay, and the Arctic Ocean.
416 However, errors stay high in the coastal regions, Northwestern passage and Hudson Bay that contribute to the high
417 total RMSEs in this region. Results are similar for the region 60°N-40°N, where correlation coefficients increased
418 from 0.3 to 0.87 on average over these zones (Fig. 10c, d). The tropical region 20°N-20°S shows correlation coefficient
419 up to 0.97 for all driver sets except DS1. However, RMSEs are high in the tropical region, about 0.11 μmol/L on
420 average (Fig. 9b), with RMSEs values of 0.2 μmol/L in the Tropical Eastern Pacific and Bay of Bengal in DS25 (Fig.
421 10b). The high RMSEs in the Tropical Eastern Pacific can indicate insufficient data in a region of high interannual
422 variability to correctly reconstruct POC_s distribution. The region of the Southern Ocean (>60°S) shows the lowest
423 correlation coefficients (in the range of 0.64-0.69) and RMSEs (in the range 0.023-0.044 μmol/L) for POC_s (Fig. 9a,
424 b). The inclusion of PFTs in the driver set significantly improves the RMSE in the region around 40°S for small
425 (POC_s) particulate organic carbon. The statistics are improved by about 75% in the region 40°S-60°S with RMSE
426 decreasing from 0.18 (DS1) to 0.03 (DS25) and the correlation coefficient increasing from 0.22 (DS1) to 0.84 (DS25),
427 on average (Fig. 9a, b; Fig. 10). The improvements in the Southern region are related to the role of zooplankton in the
428 carbon flux in this area (Le Quéré et al., 2016; Wright et al., 2021).



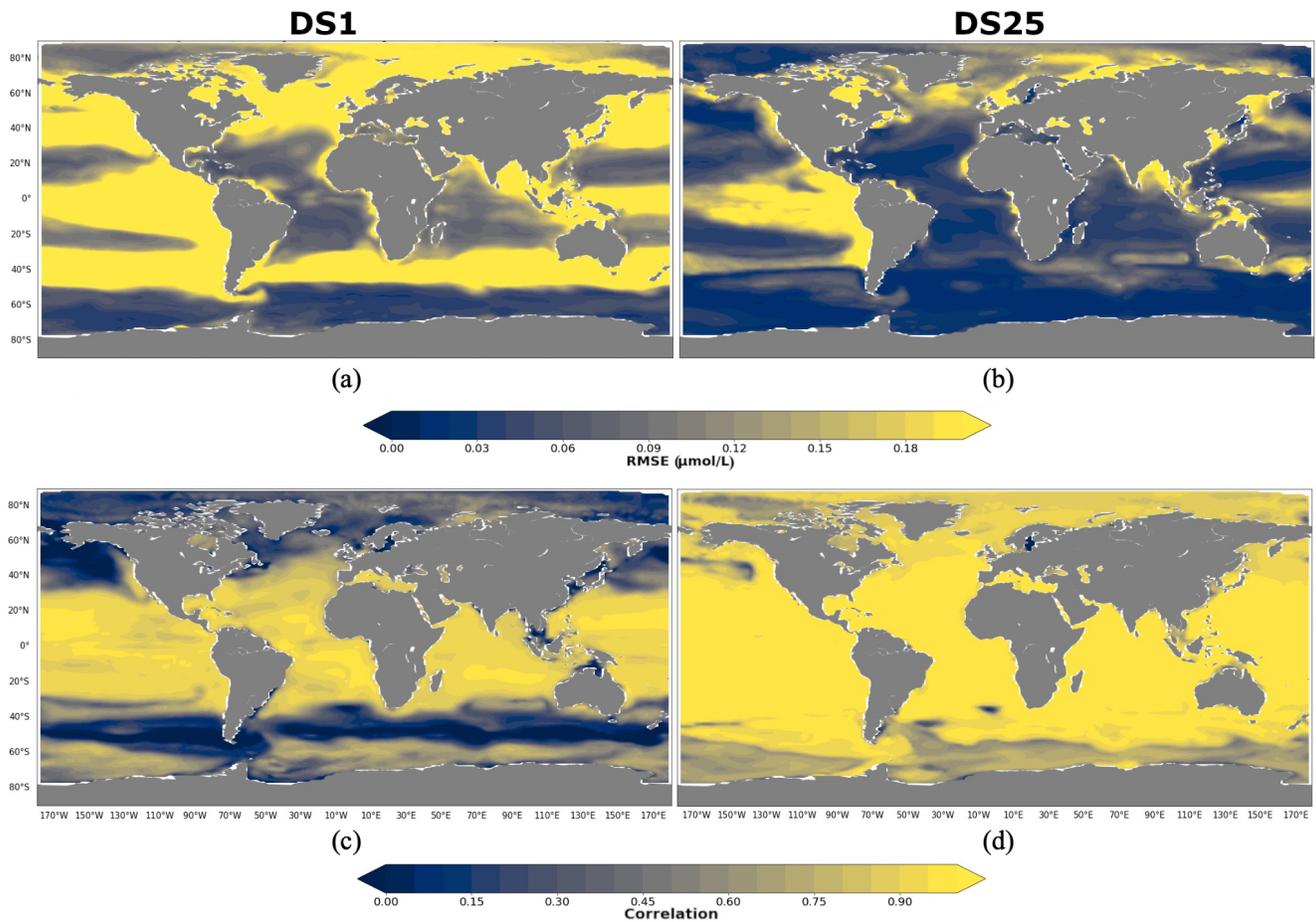
429
 430 **Figure 9. Correlations and RMSE averaged over latitude zones between PlankTOM12 and XGBoost reconstruction over the global ocean for 2009-2013: (a, c) - correlation coefficient, (b, d) - RMSE in $\mu\text{mol/l}$ (b, d); (a, b) - small particulate**
 431 **organic carbon (POCs), (c, d) - large particulate organic carbon (POC_L).**
 432

433 In POC_L reconstruction, DS1 also shows the lowest correlation coefficients (0.35-0.75) and the highest RMSEs (0.027-
 434 0.47 $\mu\text{mol/L}$) (Fig. 9c, d). DS25 shows the best results on average, with the correlation coefficient varying between
 435 0.43 (in the region 60°S-90°S) and 0.84 (in the region 20°N-20°S), and RMSE varying between 0.021 (in the region
 436 20°S-40°S) and 0.046 (in the region 90°N-60°N) $\mu\text{mol/L}$. POC_L are reconstructed better in subtropical and tropical
 437 regions compared to high latitude zones (Fig. 9c, d).

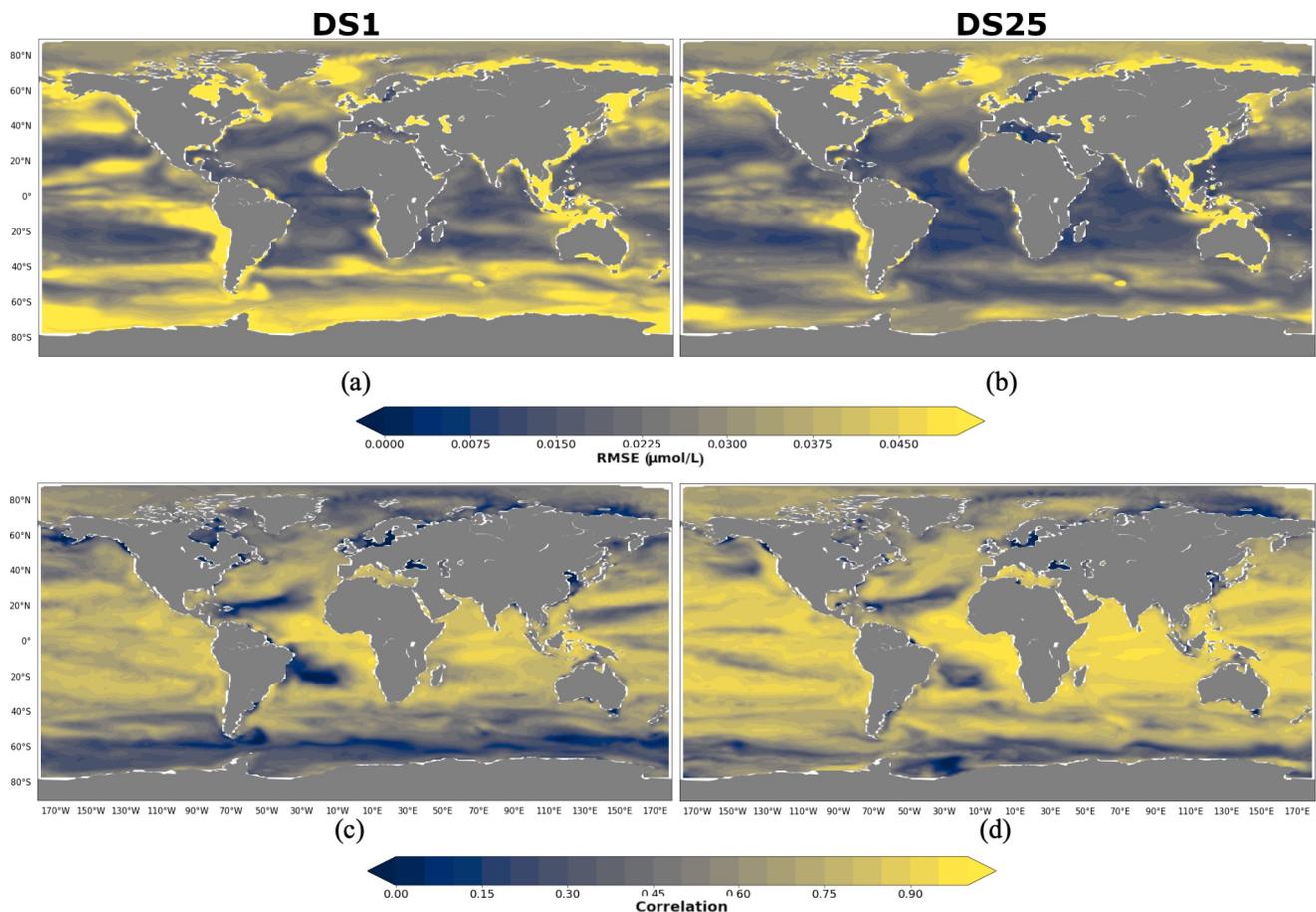
438 As for POCs, 10 driver sets (excluding DS1) show their highest RMSEs in POC_L reconstruction in the region 90°N-
 439 60°N, with values up to 0.05 $\mu\text{mol/L}$ in DS27 (Fig. 9d). Figure 11 shows maps of RMSEs (a, b) and correlation
 440 coefficients (c, d) between PlankTOM12 and reconstructed large particulate organic carbon (POC_L) by XGBoost using

441 driver sets 1 (a, c) and 25 (b, d). Contrast to POC_S reconstruction, the region 90°N-60°N does not show improvement
 442 in RMSEs for POC_L reconstruction (Fig. 11b) in DS25 compared with DS1, with still high RMSEs in Norwegian Sea,
 443 Baffin Bay, and the Arctic Ocean, and additionally for POC_L in Greenland Sea, where the algorithm did not have data
 444 for training. Similar to POC_S, errors stay high in the coastal regions, Northwestern passage and Hudson Bay that
 445 contribute to the high total RMSEs in this region.

446 Global maps of statistics suggest that the most sensible region to driver set's composition for POC_L is the Southern
 447 Ocean, as for POC_S (Fig. 11). In the 40°S-60°S region, RMSE is reduced from 0.037 μmol/L in DS1 to 0.024 μmol/L
 448 in DS25 (Fig. 9d), and the correlation coefficient is increased from 0.42 to 0.66 (Fig. 9c) on average, respectively. In
 449 the Southern region 60°S-90°S, RMSE is reduced from 0.047 μmol/L in DS1 to 0.033 μmol/L in DS25, and the
 450 correlation coefficient is increased from 0.33 to 0.42 (Fig. 9c) on average, respectively. The average correlation
 451 coefficients in this zone were found to be less than 0.5 in all tests with the highest value 0.5 in DS21. DS21 contains
 452 all PFTs and chlorophyll *a* vertical profile as drivers. The RMSE for DS21 in this region is close to the one of DS25,
 453 0.34 μmol/L and 0.33 μmol/L, respectively. It identifies the importance of chlorophyll *a* in the Southern Ocean as
 454 driver of POC_L variability.



455
 456 **Figure 10. RMSE and correlation between monthly PlankTOM12 and results of POC_S reconstruction using XGBoost over**
 457 **the period 2009-2013 for POC_S. (a, b) – RMSEs, (c, d) – correlation coefficients; (a, c) – reconstruction based on DS1**
 458 **(NoPFT); (b, d) – reconstruction based on DS25 (vertical profiles of zooplanktons, and zooplankton and phytoplankton**
 459 **averaged over MLD).**



460

461 **Figure 11. RMSE and correlation between monthly PlankTOM12 and results of POC_L reconstruction using XGBoost**
 462 **over the period 2009-2013 for POC_L . (a, b) – RMSEs, (c, d) – correlation coefficients; (a, c) – reconstruction based on**
 463 **DS1 (NoPFT); (b, d) – reconstruction based on DS25 (vertical profiles of zooplanktons, and zooplankton and**
 464 **phytoplankton averaged over MLD).**

465 The statistics of POC_S and POC_L reconstruction do not vary significantly between driver sets in all regions except in
 466 the Southern Ocean. This region is most sensitive to the composition of driver sets for both POC_S and POC_L .

467 **4. Conclusion.**

468

469 The aim of this work was to test the potential of using Machine Learning to reproduce modelled concentrations of
 470 particulate organic carbon within the ocean using the distribution of available observations. We co-localised outputs
 471 of the PlankTOM12 global biogeochemical ocean model with the positions of observations of small (POC_S) and large
 472 (POC_L) particulate organic carbon concentrations. Using PlankTOM outputs as references we could identify the best
 473 ML method for POC reconstruction and estimate method's accuracy in regions with poor observational cover.

474 We tested two ML methods to reconstruct POC_S and POC_L : the XGBoost regressor and Random Forest. Both methods
 475 are algorithms based on decision trees. XGBoost overperformed Random Forest by about 9% on average for POC_S
 476 reconstruction and by about 3% on average for POC_L reconstruction. XGBoost regressor builds the model sequentially
 477 improving it at each iterative step. At each iteration, XGBoost regressor analyses the prediction and gives more weight
 478 to the data where the fit is still wrong. It is a good tool for an unbalanced data set, like in our case where the data of
 479 particulate organic carbon concentration are sparse in time and space.

480 We tested the influence of a wide range of environmental and ecosystem drivers on POC_S and POC_L reconstruction.
 481 The introduction of Plankton Functional Types (PFTs) in the driver set greatly improves the fit and shows a linkage
 482 between surface ecosystem structure and particulate organic carbon distribution within the ocean interior. We

483 improved the accuracy of POC_s reconstruction by 59% on RMSE, 63% on absolute bias and by 52% on correlation
484 by introducing Plankton Functional Types (PFTs) in the driver sets (from the comparison of DS1 and DS25). The
485 presence of PFTs in the driver sets also improved the accuracy of POC_L reconstruction by 22% on RMSE, absolute
486 bias and correlation (from the comparison of DS1 and DS25). POC_s variability mostly depends on the depth level,
487 vertical profiles of microzooplankton, temperature and PO₄. POC_L variability depends on the depth level, MLD,
488 chlorophyll *a* averaged over MLD, vertical profiles of temperature, microzooplankton, phaeocystis and PO₄.
489 Additionally, we identified that chlorophyll *a* in driver sets improves the POC_L reconstruction in the Southern Ocean.

490 Despite the good accuracy over the global ocean on average, the statistics are worse in the coastal regions and in the
491 Tropical Eastern Pacific. The coastal regions suffer from the lack of data to represent the coastal dynamics. Therefore
492 the ML reconstructions assign open-ocean processes to coastal regions, leading to significant biases. The Tropical
493 Eastern Pacific is a region of strong interannual variability and the sparse measurements in time make it harder to
494 capture this variability correctly. Other regions with poor coverage by observations - the Eastern Indian Ocean, the
495 Western Pacific Ocean and the Southern Ocean - show the statistics of reconstruction comparable to one from regions
496 with a good cover - regions in the Atlantic Ocean. However, we found that the Southern Ocean is a more sensible
497 region to the driver set's composition. The observational data is particularly sparse in this region and our analysis
498 suggests that identifying the drivers of importance based on real dataset will be difficult.

499 Here we showed that the XGBoost regressor and Random Forest are suitable for this problem and can reconstruct
500 modelled POC_s and POC_L with appropriate accuracy. This is evidenced from the globally averaged correlation
501 coefficient up to 0.88 for POC_s and 0.68 for POC_L, and the globally averaged RMSE up to 20 % (0.08 μmol/L) of
502 standard deviation of PlankTOM12 POC_s, and 65% (0.028 μmol/L) of standard deviation of PlankTOM12 POC_L. ML
503 outputs represent well the spatial patterns of POC_s and POC_L distribution. However, the validity of the approach on
504 observations is dependent on the availability of co-located information on the drivers of importance. For some drivers
505 this should be possible (e.g. environmental conditions and chlorophyll *a*), while for other drivers information is more
506 sparse (e.g. the PFTs). Our analysis suggests that additional PFT observations would help provide broader insights
507 into the distribution of POC in the ocean. The next step of this work is to apply ML to real data using methods from
508 the present study. Testing the present ML approach on observations will also help provide suggestions for an optimal
509 set of drivers that can be measured specifically for POC reconstruction. For example, based on model results only,
510 our results suggest that microzooplankton concentration is particularly important and should be measured more
511 systematically, especially in the regions of high interannual variability. Likewise, this work provides information on
512 the variables that are less important in POC variability, like vertical profiles of gelatinous zooplankton, or mixed
513 phytoplankton for POC_s and coccolithophore for POC_L, and, thus, less important to be measured in this context. These
514 results will need to be tested with observations before firmly confirming the validity of the drivers. The validated
515 driver sets can help guide observational programs. In addition, recent advances in plankton imaging (Irisson et al.,
516 2022; Lombard et al., 2019; Orenstein et al., 2022) and omics (Faure et al., 2021) will soon provide a new global set
517 of data to estimate PFT concentrations across ocean basins allowing to better identify potential biological drivers of
518 POC variability. The new available data of PFTs will significantly facilitate the application of ML methods, such as
519 the one developed here, to observational data.

520 The relationships between key variables and surrounding conditions based on Machine Learning can provide a new
521 way for establishing parameters in ocean model parameterisation. The parameters can be time and space dependent
522 and, thus, vary from one region to another better representing the physics. Relationship between POC concentration
523 and environmental and ecosystem conditions can help to replace parameters in parameterised sinking velocity in
524 PlankTOM. The reconstructed POC concentration over the global ocean will contribute to the reconstruction of
525 porosity and opacity of particles that are key variables in the sinking matter velocity.

526 This study provides insights on the drivers that may be responsible for POC_s and POC_L variability and regional
527 dependencies. However, the dependencies are simply returning the outcome of complex ecosystem processes among
528 the drivers as represented in the PlankTOM12 model. Although these processes are based on current understanding
529 and a broad range of observations (Le Quéré et al., 2016; Wright et al., 2021; Buitenhuis et al., 2019), they remain
530 results from a model output. Observations could reveal different drivers that are important for POC_s and POC_L.
531 Depending on data availability and its time and space resolution, the final product based on observations should
532 provide new insights on the drivers that govern particulate organic carbon concentration in the real ocean.

533

534 **Data and code availability.** PlankTOM12 data used within this study are available at
535 <https://doi.org/10.5281/zenodo.7324781>. UVP5 data can be found at <https://doi.org/10.1594/PANGAEA.924375> (R.
536 Kiko et al., 2021). Codes for data preparation, development of machine learning methods and tests of different driver
537 sets as well as codes that provide figures shown in the article can be found at <https://doi.org/10.5281/zenodo.7326992>.

538 **Author contribution.** All authors contributed to the development of the methodology. ADS, CLQ, ETB designed the
539 experiments, and ADS carried them out. ADS developed codes and performed the simulations. ADS prepared the
540 paper with contributions from all coauthors.

541 **Acknowledgements.** The authors would like to thank Jean-Olivier Irisson for his contribution in the development of
542 the methodology. ADS, ETB and CLQ acknowledge support from Royal Society (grant RP\R1\191063), and NERC
543 Marine Frontiers project (grant NE/V011103/1) for CLQ. RK acknowledges support via a “Make Our Planet Great
544 Again” grant from the French National Research Agency within the “Programme d’Investissements d’Avenir” (grant
545 no. ANR-19-MPGA-0012) and by the Heisenberg program of the German Science Foundation under project number
546 469175784.

547 **References**

548 Alldredge, A.: The carbon, nitrogen and mass content of marine snow as a function of aggregate size, *Deep Sea*
549 *Research Part I: Oceanographic Research Papers*, 45, 529–541, [https://doi.org/10.1016/S0967-0637\(97\)00048-4](https://doi.org/10.1016/S0967-0637(97)00048-4),
550 1998.

551 Batten, S. D., Abu-Alhaja, R., Chiba, S., Edwards, M., Graham, G., Jyothibabu, R., Kitchener, J. A., Koubbi, P.,
552 McQuatters-Gollop, A., Muxagata, E., Ostle, C., Richardson, A. J., Robinson, K. V., Takahashi, K. T., Verheye, H.
553 M., and Wilson, W.: A Global Plankton Diversity Monitoring Program, *Front. Mar. Sci.*, 6, 321,
554 <https://doi.org/10.3389/fmars.2019.00321>, 2019.

555 Biau, G.: Analysis of Random Forest model, *Journal of Machine Learning Research*, 13(38), 1063–1095, 2012.

556 Buitenhuis, E. T., Hashioka, T., and Quéré, C. L.: Combined constraints on global ocean primary production using
557 observations and models: OCEAN PRIMARY PRODUCTION, *Global Biogeochem. Cycles*, 27, 847–858,
558 <https://doi.org/10.1002/gbc.20074>, 2013.

559 Buitenhuis, E. T., Le Quéré, C., Bednaršek, N., and Schiebel, R.: Large Contribution of Pteropods to Shallow CaCO₃
560 Export, *Global Biogeochem. Cycles*, 33, 458–468, <https://doi.org/10.1029/2018GB006110>, 2019.

561 Denvil-Sommer, A., Gehlen, M., Vrac, M., and Mejia, C.: LSCE-FFNN-v1: a two-step neural network model for the
562 reconstruction of surface ocean pCO₂ over the global ocean, *Geosci. Model Dev.*, 12, 2091–2105,
563 <https://doi.org/10.5194/gmd-12-2091-2019>, 2019.

564 Faure, E., Ayata, S.-D., and Bittner, L.: Towards omics-based predictions of planktonic functional composition from
565 environmental data, *Nature Communications*, 12(1), 4361, <https://doi.org/10.1038/s41467-021-24547-1>, 2021.

566
567 Friedlingstein, P., Jones, M. W., O’Sullivan, M., Andrew, R. M., Bakker, D. C. E., Hauck, J., Le Quéré, C., Peters, G.
568 P., Peters, W., Pongratz, J., Sitch, S., Canadell, J. G., Ciais, P., Jackson, R. B., Alin, S. R., Anthoni, P., Bates, N. R.,
569 Becker, M., Bellouin, N., Bopp, L., Chau, T. T. T., Chevallier, F., Chini, L. P., Cronin, M., Currie, K. I., Decharme,
570 B., Djutouang, L. M., Dou, X., Evans, W., Feely, R. A., Feng, L., Gasser, T., Gilfillan, D., Gkritzalis, T., Grassi,
571 G., Gregor, L., Gruber, N., Gürses, Ö., Harris, I., Houghton, R. A., Hurtt, G. C., Iida, Y., Ilyina, T., Luijkx, I. T., Jain,
572 A., Jones, S. D., Kato, E., Kennedy, D., Klein Goldewijk, K., Knauer, J., Korsbakken, J. I., Körtzinger, A.,
573 Landschützer, P., Lauvset, S. K., Lefèvre, N., Lienert, S., Liu, J., Marland, G., McGuire, P. C., Melton, J. R., Munro,
574 D. R., Nabel, J. E. M. S., Nakaoka, S.-I., Niwa, Y., Ono, T., Pierrot, D., Poulter, B., Rehder, G., Resplandy, L.,
575 Robertson, E., Rödenbeck, C., Rosan, T. M., Schwinger, J., Schwingshackl, C., Séférian, R., Sutton, A. J., Sweeney,
576 C., Tanhua, T., Tans, P. P., Tian, H., Tilbrook, B., Tubiello, F., van der Werf, G. R., Vuichard, N., Wada, C.,

577 Wanninkhof, R., Watson, A. J., Willis, D., Wiltshire, A. J., Yuan, W., Yue, C., Yue, X., Zaehle, S., and Zeng, J.:
578 Global Carbon Budget 2021, *Earth Syst. Sci. Data*, 14, 1917–2005, <https://doi.org/10.5194/essd-14-1917-2022>, 2022.

579 Friedrich, T. and Oschlies, A.: Basin-scale pCO₂ maps estimated from ARGO float data: A model study, *J. Geophys.*
580 *Res.*, 114, C10012, <https://doi.org/10.1029/2009JC005322>, 2009.

581 Gorsky, G., Aldorf, C., Kage, M., Picheral, M., Garcia, Y., and Favole, J.: Vertical distribution of suspended
582 aggregates determined by a new underwater video profiler, 1992.

583 Gorsky, G., Picheral, M., and Stemmann, L.: Use of the Underwater Video Profiler for the Study of Aggregate
584 Dynamics in the North Mediterranean, *Estuarine, Coastal and Shelf Science*, 50, 121–128,
585 <https://doi.org/10.1006/ecss.1999.0539>, 2000.

586 Guidi, L., Jackson, G. A., Stemmann, L., Miquel, J. C., Picheral, M., and Gorsky, G.: Relationship between particle
587 size distribution and flux in the mesopelagic zone, *Deep Sea Research Part I: Oceanographic Research Papers*, 55,
588 1364–1374, <https://doi.org/10.1016/j.dsr.2008.05.014>, 2008.

589 Guidi, L., Chaffron, S., Bittner, L., Eveillard, D., Larhlimi, A., Roux, S., Darzi, Y., Audic, S., Berline, L., Brum, J.
590 R., Coelho, L. P., Espinoza, J. C. I., Malviya, S., Sunagawa, S., Dimier, C., Kandels-Lewis, S., Picheral, M., Poulain,
591 J., Searson, S., Tara Oceans Consortium Coordinators, Stemmann, L., Not, F., Hingamp, P., Speich, S., Follows, M.,
592 Karp-Boss, L., Boss, E., Ogata, H., Pesant, S., Weissenbach, J., Wincker, P., Acinas, S. G., Bork, P., de Vargas, C.,
593 Iudicone, D., Sullivan, M. B., Raes, J., Karsenti, E., Bowler, C., and Gorsky, G.: Plankton networks driving carbon
594 export in the oligotrophic ocean, *Nature*, 532, 465–470, <https://doi.org/10.1038/nature16942>, 2016.

595 Hood, R. R., Laws, E. A., Armstrong, R. A., Bates, N. R., Brown, C. W., Carlson, C. A., Chai, F., Doney, S. C.,
596 Falkowski, P. G., Feely, R. A., Friedrichs, M. A. M., Landry, M. R., Keith Moore, J., Nelson, D. M., Richardson, T.
597 L., Salihoglu, B., Schartau, M., Toole, D. A., and Wiggert, J. D.: Pelagic functional group modeling: Progress,
598 challenges and prospects, *Deep Sea Research Part II: Topical Studies in Oceanography*, 53, 459–512,
599 <https://doi.org/10.1016/j.dsr2.2006.01.025>, 2006.

600 Irisson, J.-O., Ayata, S.-D., Lindsay, D. J., Karp-Boss, L., and Stemmann, L.: Machine Learning for the Study of
601 Plankton and Marine Snow from Images, *Annual Review of Marine Science*, 14(1), 277–301.
602 <https://doi.org/10.1146/annurev-marine-041921-013023>, 2022.

603

604 Kalnay, E., Kanamitsu, M., Kistler, R., Collins, W., Deaven, D., Gandin, L., Iredell, M., Saha, S., White, G., and
605 Woollen, J.: The NCEP/NCAR 40-year reanalysis project, *B. Am. Meteorol. Soc.*, 77, 437–472, 1996.

606 Kiko, R., Picheral, M., Antoine, D., Babin, M., Berline, L., Biard, T., Boss, E., Brandt, P., Carlotti, F., Christiansen,
607 S., Coppola, L., de la Cruz, L., Diamond-Riquier, E., de Madron, X. D., Elineau, A., Gorsky, G., Guidi, L., Hauss, H.,
608 Irisson, J.-O., Karp-Boss, L., Karstensen, J., Kim, D., Lekanoff, R. M., Lombard, F., Lopes, R. M., Marec, C.,
609 McDonnell, A., Niemeyer, D., Noyon, M., O'Daly, S., Ohman, M. D., Pretty, J. L., Rogge, A., Searson, S., Shibata,
610 M., Tanaka, Y., Tanhua, T., Taucher, J., Trudnowska, E., Turner, J. S., Waite, A. M., and Stemmann, L.: The global
611 marine particle size distribution dataset obtained with the Underwater Vision Profiler 5 - version 1,
612 <https://doi.org/10.1594/PANGAEA.924375>, 2021.

613 Kiko, R., Picheral, M., Antoine, D., Babin, M., Berline, L., Biard, T., Boss, E., Brandt, P., Carlotti, F., Christiansen,
614 S., Coppola, L., de la Cruz, L., Diamond-Riquier, E., Durrieu de Madron, X., Elineau, A., Gorsky, G., Guidi, L.,
615 Hauss, H., Irisson, J.-O., Karp-Boss, L., Karstensen, J., Kim, D., Lekanoff, R. M., Lombard, F., Lopes, R. M., Marec,
616 C., McDonnell, A. M. P., Niemeyer, D., Noyon, M., O'Daly, S. H., Ohman, M., Pretty, J. L., Rogge, A., Searson, S.,
617 Shibata, M., Tanaka, Y., Tanhua, T., Taucher, J., Trudnowska, E., Turner, J. S., Waite, A., and Stemmann, L.: A
618 global marine particle size distribution dataset obtained with the Underwater Vision Profiler 5, *Earth Syst. Sci. Data*

619 Discuss. [preprint], <https://doi.org/10.5194/essd-2022-51>, 14, 4315–4337, [https://doi.org/10.5194/essd-14-4315-](https://doi.org/10.5194/essd-14-4315-2022)
620 2022, 2022.

621 Kirchman, D. L.: Growth Rates of Microbes in the Oceans, *Annu. Rev. Mar. Sci.*, 8, 285–309,
622 <https://doi.org/10.1146/annurev-marine-122414-033938>, 2016.

623 Landschützer, P., Gruber, N., Bakker, D. C. E., Schuster, U., Nakaoka, S., Payne, M. R., Sasse, T. P., and Zeng, J.: A
624 neural network-based estimate of the seasonal to inter-annual variability of the Atlantic Ocean carbon sink,
625 *Biogeosciences*, 10, 7793–7815, <https://doi.org/10.5194/bg-10-7793-2013>, 2013.

626 Le Quéré, C., Harrison, S. P., Colin Prentice, I., Buitenhuis, E. T., Aumont, O., Bopp, L., Claustre, H., Cotrim Da
627 Cunha, L., Geider, R., Giraud, X., Klaas, C., Kohfeld, K. E., Legendre, L., Manizza, M., Platt, T., Rivkin, R. B.,
628 Sathyendranath, S., Uitz, J., Watson, A. J., and Wolf-Gladrow, D.: Ecosystem dynamics based on plankton functional
629 types for global ocean biogeochemistry models, *Global Change Biology*, 0, 2016–2040,
630 <https://doi.org/10.1111/j.1365-2486.2005.1004.x>, 2005.

631 Le Quéré, C., Buitenhuis, E. T., Moriarty, R., Alvain, S., Aumont, O., Bopp, L., Chollet, S., Enright, C., Franklin, D.
632 J., Geider, R. J., Harrison, S. P., Hirst, A. G., Larsen, S., Legendre, L., Platt, T., Prentice, I. C., Rivkin, R. B., Saille, S.,
633 Sathyendranath, S., Stephens, N., Vogt, M., and Vallina, S. M.: Role of zooplankton dynamics for Southern Ocean
634 phytoplankton biomass and global biogeochemical cycles, *Biogeosciences*, 13, 4111–4133,
635 <https://doi.org/10.5194/bg-13-4111-2016>, 2016.

636 Lombard, F., Boss, E., Waite, A. M., Vogt, M., Uitz, J., Stemmann, L., Sosik, H. M., Schulz, J., Romagnan, J.-B.,
637 Picheral, M., Pearlman, J., Ohman, M. D., Niehoff, B., Möller, K. O., Miloslavich, P., Lara-Lpez, A., Kudela, R.,
638 Lopes, R. M., Kiko, R., Karp-Boss, L., Jaffe, J. S., Iversen, M. H., Irisson, J.-O., Fennel, K., Hauss, H., Guidi, L.,
639 Gorsky, G., Giering, S. L. C., Gaube, P., Gallagher, S., Dubelaar, G., Cowen, R. K., Carlotti, F., Briseño-Avena, C.,
640 Berline, L., Benoit-Bird, K., Bax, N., Batten, S., Ayata, S. D., Artigas, L. F., and Appeltans, W.: Globally Consistent
641 Quantitative Observations of Planktonic Ecosystems, *Front. Mar. Sci.*, 6, 196,
642 <https://doi.org/10.3389/fmars.2019.00196>, 2019.

643 Mutshinda, C., Finkel, Z., Widdicombe, C., and Irwin, A.: Phytoplankton traits from long-term oceanographic time-
644 series, *Mar. Ecol. Prog. Ser.*, 576, 11–25, <https://doi.org/10.3354/meps12220>, 2017.

645 Orenstein, E. C., Ayata, S., Maps, F., Becker, É. C., Benedetti, F., Biard, T., et al.: Machine learning techniques to
646 characterize functional traits of plankton from image data. *Limnology and Oceanography*, 67(8), 1647–1669.
647 <https://doi.org/10.1002/lno.12101>, 2022.
648

649 Picheral, M., Guidi, L., Stemmann, L., Karl, D. M., Iddaoud, G., and Gorsky, G.: The Underwater Vision Profiler 5:
650 An advanced instrument for high spatial resolution studies of particle size spectra and zooplankton: Underwater vision
651 profiler, *Limnol. Oceanogr. Methods*, 8, 462–473, <https://doi.org/10.4319/lom.2010.8.462>, 2010.

652 Sauzède, R., Claustre, H., Uitz, J., Jamet, C., Dall’Olmo, G., D’Ortenzio, F., Gentili, B., Poteau, A., and Schmechtig,
653 C.: A neural network-based method for merging ocean color and Argo data to extend surface bio-optical properties to
654 depth: Retrieval of the particulate backscattering coefficient, *J. Geophys. Res. Oceans*, 121, 2552–2571,
655 <https://doi.org/10.1002/2015JC011408>, 2016.

656 Sauzède, R., Bittig, H. C., Claustre, H., Pasqueron de Fommervault, O., Gattuso, J.-P., Legendre, L., and Johnson, K.
657 S.: Estimates of Water-Column Nutrient Concentrations and Carbonate System Parameters in the Global Ocean: A
658 Novel Approach Based on Neural Networks, *Front. Mar. Sci.*, 4, 128, <https://doi.org/10.3389/fmars.2017.00128>, 2017.

659 Sauzède, R., Johnson, J. E., Claustre, H., Camps-Valls, G., and Ruescas, A. B.: ESTIMATION OF OCEANIC
660 PARTICULATE ORGANIC CARBON WITH MACHINE LEARNING, *ISPRS Ann. Photogramm. Remote Sens.*
661 *Spatial Inf. Sci.*, V-2–2020, 949–956, <https://doi.org/10.5194/isprs-annals-V-2-2020-949-2020>, 2020.

662 Schlitzer, R.: Carbon export fluxes in the Southern Ocean: results from inverse modeling and comparison with
663 satellite-based estimates, *Deep Sea Research Part II: Topical Studies in Oceanography*, 49, 1623–1644,
664 [https://doi.org/10.1016/S0967-0645\(02\)00004-8](https://doi.org/10.1016/S0967-0645(02)00004-8), 2002.

665 Sunagawa, S., Acinas, S. G., Bork, P., Bowler, C., Tara Oceans Coordinators, Acinas, S. G., Babin, M., Bork, P.,
666 Boss, E., Bowler, C., Cochrane, G., de Vargas, C., Follows, M., Gorsky, G., Grimsley, N., Guidi, L., Hingamp, P.,
667 Iudicone, D., Jaillon, O., Kandels, S., Karp-Boss, L., Karsenti, E., Lescot, M., Not, F., Ogata, H., Pesant, S., Poulton,
668 N., Raes, J., Sardet, C., Sieracki, M., Speich, S., Stemmann, L., Sullivan, M. B., Sunagawa, S., Wincker, P., Eveillard,
669 D., Gorsky, G., Guidi, L., Iudicone, D., Karsenti, E., Lombard, F., Ogata, H., Pesant, S., Sullivan, M. B., Wincker, P.,
670 and de Vargas, C.: Tara Oceans: towards global ocean ecosystems biology, *Nat Rev Microbiol*, 18, 428–445,
671 <https://doi.org/10.1038/s41579-020-0364-5>, 2020.

672 Telszewski, M., Chazottes, A., Schuster, U., Watson, A. J., Moulin, C., Bakker, D. C. E., Gonzalez-Davila, M.,
673 Johannessen, T., Kortzinger, A., Santana-Casiano, M., Wallace, D. W. R., and Wanninkhof, R.: Estimating the
674 monthly pCO₂ distribution in the North Atlantic using a self-organizing neural network, 17, 2009.

675 Wright, R. M., Le Quéré, C., Buitenhuis, E., Pitois, S., and Gibbons, M. J.: Role of jellyfish in the plankton ecosystem
676 revealed using a global ocean biogeochemical model, *Biogeosciences*, 18, 1291–1320, [https://doi.org/10.5194/bg-18-](https://doi.org/10.5194/bg-18-1291-2021)
677 1291-2021, 2021.